

CompSci 590.01 Spring 2023

Causal Inference in Data Analysis with Applications to Fairness and Explanations

> Lecture 7: Explanations & Causal Explanations

> > Sudeepa Roy

Reading

• Galhotra et al. et al., SIGMOD (2021): Explaining Black-Box Algorithms Using Probabilistic Contrastive Counterfactuals (<u>https://dl.acm.org/doi/10.1145/3448016.3458455</u>)

Acknowledgement (big thanks!):

Slides are from a tutorial by Prof. Babak Salimi (UCSD), Prof. Romila Pradhan (Purdue), Dr. Sainyam Galhotra (U. Chicago), and graduate student Aditya Lahiri with small modifications from a tutorial they gave at ICDE'22 and SIGMOD'22, and from a joint short course at the Reasoning Web'22 summer school with Prof. Salimi

Announcements: Feb 9

- Please see presentation guideline on Ed
- Initial project feedback coming soon
- Send me an email by today with your interests if you are currently only yourself in a project team and looking for 1-2 teammates

Causal Inference and Explanations

Explanations – another important topic in Responsible Data Science

PURDUE



THE UNIVERSITY OF CHICAGO

Explainable AI: Foundations, Applications, and Opportunities for Data Management Research











Romila Pradhan Assistant Professor Purdue

Aditya Lahiri Graduate Student UCSD

Sainyam Galhotra Postdoc UChicago

Babak Salimi Assistant Professor UCSD

https://explainable-ai-tutorial.github.io/ [SIGMOD'22, ICDE'22]

We interact with algorithmic decision-making on daily basis

Sophisticated ML models shown to be highly accurate for many applications



- Complex models → difficult to trace decisions back to questions about why and how they were made
- ML systems are often opaque– complexity, proprietary

Human-understandable explanations of outcomes of algorithmic decision-making systems

 A powerful tool for answering How? and Why? questions about algorithmic systems – also for legal & ethical concerns

Not a new topic!

HPP-79-11 (Working Paper)

March 1979

Transfer of Expertise: A Theme for Al Research

A key idea in our current approach to building expert systems is that these programs should not only be able to apply the corpus of expert knowledge to specific problems, but they should also be able to interact with the users and experts just as humans do when they learn, explain, and teach what they know. We will show, as we review the major

Explanation in Second Generation Expert Systems

William R. Swartout¹ and Johanna D. Moore²

¹ USC/Information Sciences Institute 4676 Admiralty Way Marina del Rey, CA 90292

² University of Pittsburgh Computer Science Department and Learning Research and Development Center Pittsburgh, PA 15260

Abstract. What is needed for good explanation? This paper begins by considering some desiderata for expert system explanation. These desiderata concern not only the form and content of the explanations, but also the impact of explanation generation on the expert system itself how it is built and how it performs. In this paper, we use these desiderata as a yardstick for measuring progress in the field. The paper describes

A Theory of Diagnosis from First Principles



Raymond Reiter

Department of Computer Science, University of Toronto, Toronto, Ontario, Canada M5S 1A4; The Canadian Institute for Advanced Research

Recommended by Johan de Kleer and Daniel G. Bobrow

ABSTRACT

Suppose one is given a description of a system, together with an observation of the system's behaviour which conflicts with the way the system is meant to behave. The diagnostic problem is to determine those components of the system which, when assumed to be functioning abnormally, will explain the discrepancy between the observed and correct system behaviour.

we propose a general theory for this problem. The theory requires only that the system be described in a suitable logic. Moreover, there are many such suitable logics, e.g. first-order, temporal, dynamic, etc. As a result, the theory accommodates diagnostic reasoning in a wide variety of practical settings, including digital and analogue circuits, medicine, and database updates. The theory leads to an algorithm for computing all diagnoses, and to various results concerning principles of measurement for discriminating among competing diagnoses. Finally, the theory reveals close connections between diagnostic reasoning and nonmonotonic reasoning.

as a yardstick for measuring progress in the field. The paper describes



Applications of XAI

Overcoming Racial Bias In AI Systems And Startlingly Even In AI Self-Driving Cars Racial bias in a medical algorithm favors white patients over sicker black patients

AI expert calls for end to UK use of 'racially biased' algorithms

AI Bias Could Put Women's Lives At Risk - A Challenge For Regulators

Gender bias in Al: building fairer algorithms

Bias in AI: A problem recognized but still unresolved

Amazon, Apple, Google, IBM, and Microsoft worse at transcribing black people's voices than white people's with Al voice recognition, study finds

Millions of black people affected by racial bias in health-care algorithms

Study reveals rampant racism in decision-making software used by US hospitals – and highlights ways to correct it.

When It Comes to Gorillas, Google Photos Remains Blind

Google promised a fix after its photo-categorization software labeled black people as gorillas in 2015. More than two years later, it hasn't found one.

The Week in Tech: Algorithmic Bias Is Bad. Uncovering It Is Good.

Google 'fixed' its racist algorithm by removing gorillas from its image-labeling tech

Artificial Intelligence has a gender bias problem – just ask Siri

The Best Algorithms Struggle to Recognize Black Faces Equally

US government tests find even top-performing facial recognition systems misidentify blacks at rates five to 10 times higher than they do whites.

Image source: https://towardsdatascience.com/algorithmbias-in-artificial-intelligence-needs-to-be-discussed-and-

Applications of XAI

Accuracy, trust, recourse and compliance with the law



Machine Learning In Healthcare Industry

Companies should commit to ensuring systems that could fall under GDPR, including AI, will be compliant. The threat of sizeable fines of €20 million or 4% of global turnover provides a sharp incentive.

"

Article 22 of GDPR empowers individuals with the right to demand an explanation of how an AI system made a decision that affects them.



,,

"The data subject shall have the **right** not to be subject to a decision based solely on automated processing..."

"... monitor city use of algorithmic decision-making and provide recommendations..."

Automated Decision Systems Task Force

required companies to study algorithms they use, identify bias in these systems and fix any discrimination or bias they find



Algorithmic Accountability Act 2019



Applications of XAI

Debugging ML model





eXplainable AI approaches for debugging and diagnosis.

Workshop @ NeurIPS2021 | 14 December

Schedule

12

How are such decisions explained?



H. Lakkaraju, J. Adebayo, S. Singh. NeurIPS 2020.

Feature-based explanations





Data-based explanations

- Feature-based explanations fall short in generating diagnostic explanations
- Data-based explanations: parts of the training data are responsible for an unexpected and discriminatory behavior of an ML model?



Causal Explanations

Counterfactual explanations



Counterfactual Explanations

- Goal: Minimum change in input attributes to flip a model's prediction
- Popular Techniques
 - Feature-based explanations: Feature score captures its likelihood to change outcome
 - Instance-level explanations: Identify counterfactual scenarios to change outcome
 - Inverse Classification
 - Nearest-counterfactual explanation

Feature-based explanation

 Counterfactual score of a feature: Expected difference in outcome when the feature takes a value of any other individual value



Nam	e	Age	Saving	Month	 Credit Amount	Black-
Maev	/e	<25	<100 DM	24	 Existing paid duly	box ML model

Name	Age	Saving	Month	 Credit Amount	Black-
Tom	<25	500 DM	24	 Existing paid duly	box ML model
Name	Age	Saving	Month	 Credit Amount	Black-
Tim	<25	200 DM	24	 Existing paid dulv	box ML model

Inverse Classification

- Identify a close neighbor of the point which has the opposite outcome
- Growing sphere algorithm



Inverse Classification for Comparison-based Interpretability in Machine Learning Laugel et al., <u>https://arxiv.org/pdf/1712.08443.pdf</u> - figure from the paper

Limitations

- Focus on correlation between input/output
- Fail to account for the causal interaction between attributes
- Perturbations not translatable to real-world interventions
- Fall short in generating diagnostic explanations



What kind of explanations do humans seek?

"Explanatory relevant information is information that is potentially **relevant to manipulation and control**"

-James Woodward, Philosopher

The key insight is to recognize that one does not explain events per se, but that one explains why the puzzling event occurred in the target cases but not in some counterfactual contrast case.

-Denis Hilton, Psychologist

"An explanation is an **assignment of causal responsibility**" —Josephson and Josephson, Computer Scientists

"To explain an event is to provide some **information about its causal history**." "... think of a cause as something that makes a difference... **Had it been absent, its effects** – some of them, at least, and usually all – **would have been absent as well**" —David Lewis, Philosopher

CONTRASTIVE

"...the only way to ensure that the recommended change is even possible..." "...and to account for dependencies between features is to model the outcome of interest using features that directly figure into the causal mechanism" —Barocas et al. FAT 2020

Miller, Tim. "Explanation in artificial intelligence: Insights from the social sciences." *Artificial intelligence* 267 (2019): 1-38.

CAUSAL

Recap: Causal Graph



Fully specified model (functions and values of exogenous variables known)

Age = U_{Age} Gender = U_{Gender} Month = 0.4 Age + 0.8 Gender + U_{Month} Decision = Month + 0.5 Age + $U_{Decision}$

Recap: Counterfactuals/Interventions



Equivalent to removing incoming edges of credit amount and conditioning on it

$$\Pr(y_{\mathbf{X}\leftarrow\mathbf{x}} \mid \mathbf{k}) = \sum_{\mathbf{u}} \Pr(y_{\mathbf{X}\leftarrow\mathbf{x}}(\mathbf{u})) \Pr(\mathbf{u} \mid \mathbf{k})$$

Pr(u) induces a probability distribution over exogenous variables

Recap: Counterfactuals/Interventions



Pearl's three step procedure

- 1. Abduction: Calculate Pr(u|k)
- 2. Action: modify the causal graph to reflect intervention
- 3. Prediction: Calculate the probability of Y=y

Causal Explanations

- Advantages
 - Translatable to real-world interventions
 - Easily understandable
 - Robust against spurious correlations in the data
- Challenges
 - Estimation requires additional information
 - Domain knowledge
 - Causal graph

Galhotra-Pradhan-Salimi, SIGMOD'21

Lewis: Unifying Contrastive and feature attribution methods



Based on probabilistic contrastive counterfactuals

For individual(s) with attribute(s) <actual-value> for whom an algorithm made the decision <actual-outcome>, the decision would have been <foil-outcome> with probability <score> had the attribute been <counterfactual-value>.

To what extent is an attribute sufficient?



Rosa



What is the probability that loan would be approved if month had been 36 months instead of 24 months?



To what extent is an attribute necessary?





Irrfan, a bank customer Irrfan's data

What is the probability that loan would be **rejected** if month had been 24 months instead of 36 months?



48%

Irrfan's

credibility

Our explanation scores

attribute X has value x positive model outcome instances in context

Necessity score

 $\operatorname{Nec}_{x}^{x'}(\mathbf{k}) = \Pr(o'_{X \leftarrow x'} \mid x, o, \mathbf{k})$

Probability that outcome would have been negative had X been x'

Probability that Loan=Rejected if Month were 24 for instances in context for whom Loan=Approved when Month=36

Our explanation scores

$$Nec_{x}^{x'}(\mathbf{k}) = Pr(o'_{X \leftarrow x'} \mid x, o, \mathbf{k})$$

attribute X has value x'
negative model outcome
instances in context
Sufficiency
score
$$SuF_{x}^{x'}(\mathbf{k}) = Pr(o_{X \leftarrow x} \mid x', o', \mathbf{k})$$

Probability that outcome would

have been positive had X been x

Probability that Loan=Approved if Month were 36 for instances in context for whom Loan=Rejected when Month=24

Lewis

For context k

$\operatorname{Nec}_{x}^{x'}(\mathbf{k})$	$\Pr(o'_{X \leftarrow x'} \mid x, o, \mathbf{k})$	Probability that Loan=Rejected were Month=24 for whom Loan=Approved when Month=36
$\operatorname{Suf}_{x}^{x'}(\mathbf{k})$	$\Pr(o_{X \leftarrow x} \mid x', o', \mathbf{k})$	Probability that Loan=Approved were Month=36 for whom Loan=Rejected when Month=24
$\operatorname{NeSuf}_{x}^{x'}(\mathbf{k})$	$\Pr(o_{X \leftarrow x}, o'_{X \leftarrow x'} \mid \mathbf{k})$	Probability that Loan=Approved were Month=36 and Loan=Rejected were Month=24

Our explanation scores

 $\operatorname{Nec}_{x}^{x'}(\mathbf{k}) = \Pr(o'_{X \leftarrow x'} \mid x, o, \mathbf{k})$

$$\operatorname{Sur}_{x}^{x'}(\mathbf{k}) = \Pr(o_{X \leftarrow x} \mid x', o', \mathbf{k})$$

$$NeSur_{x}^{x'}(\mathbf{k}) = Pr(o_{X \leftarrow x}, o_{X \leftarrow x'}' | \mathbf{k})$$

$$NeSur_{x}^{x'}(\mathbf{k}) \leq Pr(o, \mathbf{x} | \mathbf{k}) \underbrace{Nec_{\mathbf{x}}(\mathbf{k})}_{+ Pr(o', \mathbf{x}' | \mathbf{k}) \underbrace{Sur_{\mathbf{x}}(\mathbf{k})}_{+ 1 - Pr(\mathbf{x} | \mathbf{k}) - Pr(\mathbf{x}' | \mathbf{k})}$$

NeSuf can be seen to be a linear combination of Nec and Suf scores, And for binary X, this inequality becomes an equality and NeSuf can be seen as the weighted average of Nec and Suf scores.

From scores to explanations



Global explanations



 $\operatorname{SuF}_{x}^{x'}(\mathbf{k}=\phi)$

- Scores express the global influence of attributes on algorithm's decision
- Maximum score over all pairs of an attribute's values

Local explanations



Contextual explanations



Algorithmic Recourse

- Input: Individual with negative outcome
- Goal: Identify intervention with smallest cost that can flip the outcome in the future



Name	Age	Saving	Month	 Credit Amount	Ĺ	Black-	
Maeve	<25	<100 DM	24	 Existing paid duly		model	

What should I do to flip the outcome?

Non-causal methods

- Model recourse as an optimization problem
- Linear Program:
 - Objective: Minimize the change in attributes
 - Constraints: Flip model outcome
- Assumption: All attributes are independent and do not impact each other

Causal Techniques

- Formulate as an optimization problem
- Goal: Minimize the cost of actions
- Constraints:
 - Classifier outcome is flipped
 - Modified input is causally consistent with the original input
- Assumption: Requires fully specified causal model

Karimi et al., Algorithmic Recourse: from Counterfactual Explanations to Interventions ACM FAacT'20 Lewis: How can I improve my chances of getting a loan?

User provides a set of actionable attributes A

Name	Age	Saving	Month	 Credit History
Maeve	<25	<100 DM	24	 Existing paid duly

Counterfactual recourse as interventions over actionable attributes

$$\begin{array}{c|c} \operatorname{argmin}_{\mathbf{a}\in\operatorname{Dom}(\mathbf{A})} & \operatorname{Cost}\left(\mathbf{a},\hat{\mathbf{a}}\right) \\ & \text{s.t.} & \operatorname{SuF}_{\hat{\mathbf{a}}}(\operatorname{Maeve's\ data}) \geq 85\% \end{array} \right)$$

Recommended Recourse:

Actionable Attributes	Current Value	Required Value
Credit amount	1,275 DM	3,000 – 5,000 DM
Savings	< 100 DM	500 – 1,000 DM

Can operate at sub-population level Requires causal structure but not dependence equations



Computational Challenges

NEC_x^{x'}(k)
$$\Pr(o'_{X \leftarrow x'} \mid x, o, k)$$
 $SUF_x^{x'}(k)$ $\Pr(o_{X \leftarrow x} \mid x', o', k)$ NESUF_x^{x'}(k) $\Pr(o_{X \leftarrow x}, o'_{X \leftarrow x'} \mid k)$

Requires the full structural causal model





Experiments



Causal Shapley

- SHAP will be covered by student presentations causal version
- Proposed a causal notion of the value function
- Captures direct and indirect effects of an attribute

Common Challenges & Takeaways

- Dependence on causal structure
- Fragile with respect to noise in data
 - Distribution shift
- Need to adapt to evolving needs
- Causal explanations are translatable to real-world interventions
- Challenges
 - Require additional information
 - Computationally difficult
- Recourse
 - Helps individuals to reverse algorithm's decisions
 - Generally solved as an optimization problem

Data-based explanations

Explaining model outcomes





Data-based explanations



Name	Age	Education	Marital	 Race	Gender	Hours	Income
Nazia	36	Bachelors	Unmarried	 Asian	female	40	<50K
Matt	26	Bachelors	Married	 White	male	40	≥50K
Yeji	50	Masters	Married	 Asian	male	16	≥50K
Neel	45	Masters	Unmarried	 Black	male	28	<50K

Idea: Select training data points to explain ML model behavior



Approaches for Data-based Explanations

• **Objective:** debug data as a way to debug ML models, understand and explain model behavior and model predictions

• Popular approaches:

- Data valuation
- Influence functions
- Representer points
- Prototypes and criticisms

How valuable or influential is a data point?

• How much does each training data point contribute to the learned



	Name	Age	Educa	ation	Marital		Race	Ger	nder	Hours		
	Rosa 34 Bach			lors Unmarried			Black	ferr	nale	40	1	
Rosa's data												
					-	M	odel					
Name	e Age	Edu	cation	Marit	al	Rac	odel e G	ender	Hours	s Inco	ome	1
<mark>Name</mark> Nazia	<mark>e Age</mark> a 36	Edu Bac	cation helors	Marit Unmari	<mark>al</mark> ried	Rac Asia	odel e G n fe	ender emale	Hours 40	s Inco	ome 50K	
<mark>Name</mark> Nazia Matt	e Age a 36	Edu Bac Bac	cation helors helors	Marit Unmari Marrie	<mark>al</mark> ried ed	Rac Asia Whit	odel e G n fe e I	ender emale male	Hours 40 40	<mark>s Inco</mark> <5 ≥5	ome 50K 50K	
<mark>Name</mark> Nazia Matt Yeji	e Age a 36 26 50	Edu Bac Bac Ma	cation helors helors	Marita Unmari Marrie Marrie	<mark>al</mark> ried ed	Rac Asia Whit Asia	odel e G n fe e r	ender emale male male	Hours 40 40 16	s Inco <5 ≥5 ≥5	<mark>ome</mark> 50K 50K	-
<mark>Name</mark> Nazia Matt Yeji Neel	e Age a 36 26 50 45	Edu Bac Bac Ma	cation helors helors isters	Marita Unmari Marrie Marrie Unmari	<mark>al</mark> ried ed ed	Rac Asia Whit Asia Blac	odel e G n fe e r k r	ender emale male male male	Hours 40 40 16 28	s Inco <5 ≥5 ≥5 <5	<mark>ome</mark> 50K 50K 50K	-

How valuable or influential is a data point?

• Leave-one-out approach



Cook, R. D. Detection of influential observation in linear regression. Technometrics, 19(1):15–18, 1977

How valuable or influential is a data point?

Leave-one-out approach

Original model Updated model performance performance

	Name	Age	Education	Marital	 Race	Gender	Hours	Income		7
Value (Nazia	36	Bachelors	Unmarried	 Asian	female	40	<50K) = 0.83 -	- 0.75
	Matt	26	Bachelors	Married	 White	male	40	≥50K	= 0.08	
	Yeji	50	Masters	Married	 Asian	male	16	≥50K		
	Neel	45	Masters	Unmarried	 Black	male	28	<50K		

- Might be expensive to compute for each data point
- Not reasonable when training data has duplicates

When is a data valuation approach good?

 Null Element: If adding training data point t₁ to <u>any</u> subset of training data never changes the learned model's performance:



Education Marital Gender Name Age Race Hours . . . Rosa 34 **Bachelors** Unmarried Black female 40 . . . Rosa's data 83% Probability that Rosa's income < 50KOriginal Model Education Marital Name Age Race Gender Hours Income . . . 36 **Bachelors** Unmarried <50K Nazia Asian female 40 . . . 26 **Bachelors** Married White male ≥50K Matt 40 . . . Married Yeji 50 Masters Asian ≥50K male 16 . . . 45 Masters Unmarried Black 28 <50K Neel male .

55

value(t_1) = 0

Ghorbani, A., Zou, J.. Data Shapley: Equitable Valuation of Data for Machine Learning. ICML. 2019

1. Null Element: If adding training data point t_1 to <u>any</u> subset of training data never changes the learned model's performance:



1. Null Element: If adding training data point t_1 to <u>any</u> subset of training data never changes the learned model's performance:



Symmetry: If adding training data point t_1 or t_2 to any subset of 2. training data always results in the same change in performance:



 $value(t_1) = value(t_2)$



2. **Symmetry:** If adding training data point t_1 or t_2 to any subset of training data always results in the same change in performance:



2. Symmetry: If adding training data point t_1 or t_2 to any subset of training data always results in the same change in performance:



3. Linearity: In ML, a performance metric (e.g., accuracy) is often the sum of the performance metric on individual test points

In this case, value of a datum should be sum of its value for each prediction

	Name	Age	Education	Marital	 Race	Gender	Hours	Income	
	Nazia	36	Bachelors	Unmarried	 Asian	female	40	<50K	
	Matt	26	Bachelors	Married	 White	male	40	≥50K	
	Yeji	50	Masters	Married	 Asian	male	16	≥50K	
Value (Neel	45	Masters	Unmarried	 Black	male	28	<50K) = Value for 🕯
					 				ا و
			-	-	-				+ Value for

test₁ test₂

Other approaches: Data Shapley and Influence

 Given training data D and performance metric P, the value of a data point t₁ is given by Data Shapley Value Marginal contribution of t₁



Idea of *influence* from robust statistics How do the model parameters change as we upweight z_{train} by an infinitesimal amount ϵ ? Also group influence

Ghorbani, A., Zou, J.. Data Shapley: Equitable Valuation of Data for Machine Learning. ICML. 2019

Koh, P.W. & Liang, P. Understanding Black-box Predictions via Influence Functions. ICML 2017

Basu, S., You, X., and Feizi, S. On Second-Order Group Influence Functions for Black-Box Predictions. ICML 2020

62